

CLAIMS

1 1. In a cluster of computing nodes having shared access
2 to one or more volumes of data storage using a parallel
3 file system, a method for managing the data storage,
4 comprising:

5 selecting a first one of the nodes to serve as a
6 session manager node;

7 selecting a second one of the nodes to serve as a
8 session node for a data management application;

9 creating a session of the data management
10 application on the session node by sending a message from
11 the session node to the session manager node, causing the
12 session manager node to distribute information regarding
13 the session among the nodes in the cluster; and

14 responsive to the information distributed by the
15 session manager node, receiving events at the session
16 node from the nodes in the cluster for processing by the
17 data management application.

1 2. A method according to claim 1, and comprising
2 storing the information regarding the session at the
3 session manager node.

1 3. A method according to claim 2, wherein the
2 information regarding the session is stored at both the
3 session node and at the session manager node, and
4 comprising, following a failure at the session node,
5 receiving the stored information from the session manager
6 node in order to recover the session.

1 4. A method according to claim 3, wherein at least a
2 portion of the information regarding the session is
3 stored at one or more of the nodes in the cluster other

than at the session node and the session manager node, and comprising, following a failure at the first one of the nodes, selecting a third one of the nodes to serve as the session manager node, and collecting the information from at least one of the session node and the other nodes in the cluster at which the information is stored for use by the third one of the nodes in serving as the session manager node.

5. A method according to claim 1, wherein creating the session comprises creating the session in accordance with a data management application programming interface (DMAPI) of the parallel file system, and wherein sending the message comprises invoking a session management function of the DMAPI on the session manager node.

6. A method according to claim 1, wherein the information comprises a list of the events in each file system of relevance to the data management application and respective dispositions of the events on the list, and wherein receiving the events at the session node comprises receiving messages reporting the events appearing on the list responsive to the dispositions.

7. A method according to claim 6, and comprising receiving a data management application programming interface (DMAPI) function call from one or more of the nodes other than the session node setting one or more of the dispositions.

8. A method according to claim 6, wherein the session is one of a plurality of sessions in the cluster, and wherein the session manager node coordinates a consistent partitioning of the dispositions among the sessions.

1 9. A method according to claim 1, wherein selecting the
2 second one of the nodes comprises selecting a plurality
3 of the nodes to serve as respective session nodes in a
4 plurality of data management sessions, and wherein
5 creating the session comprises informing the session
6 manager node of the plurality of the sessions, causing
7 the session manager node to distribute the information
8 regarding the plurality of the sessions.

1 10. A method according to claim 9, wherein the first one
2 of the nodes serves as one of the session nodes, in
3 addition to serving as the session manager node.

1 11. A method according to claim 9, wherein selecting the
2 plurality of the nodes comprises selecting multiple
3 session nodes for a distributed data management
4 application running in the cluster.

1 12. A method according to claim 9, wherein selecting the
2 plurality of the nodes comprises selecting the second one
3 of the nodes to serve as the respective session node for
4 a first data management application, and selecting a
5 third one of the nodes to serve as the respective session
6 node for a second data management application.

1 13. A method according to claim 1, and comprising
2 modifying the session by sending a further message from
3 the session node to the session manager node, causing the
4 session manager node to distribute a notification
5 regarding the modified session to the nodes in the
6 cluster.

1 14. A method according to claim 1, and comprising
2 destroying the session by sending a further message from
3 the session node to the session manager node, causing the

4 session manager node to distribute a notification
5 regarding the destroyed session to the nodes in the
6 cluster.

1 15. A method according to claim 1, wherein creating the
2 session of the data management application comprises
3 initiating a data migration application, so as to free
4 storage space on at least one of the volumes of data
5 storage.

1 16. A method according to claim 1, wherein the
2 information comprises a session identifier, generated at
3 the session manager node, which is unique within the
4 cluster.

1 17. In a cluster of a plurality of computing nodes
2 having shared access to one or more volumes of data
3 storage using a parallel file system, a method for
4 managing the data storage, comprising:

5 initiating sessions of a parallel data management
6 application on the plurality of the nodes, so that an
7 instance of the data management application runs on each
8 of the nodes;

9 generating a data management event responsive to a
10 request submitted to the parallel file system on at least
11 one of the nodes to perform a file operation on a file in
12 one of the volumes of data storage;

13 handling the event by means of the instance of the
14 data management application running on the at least one
15 of the nodes.

1 18. A method according to claim 17, and comprising
2 sending an event message from the at least one of the
3 nodes to the other nodes, so as to inform the data

4 management application sessions on the other nodes of the
5 event.

1 19. A method according to claim 17, wherein generating
2 the data management event comprises running a user
3 application on the at least one of the nodes, and
4 receiving the request from the user application.

1 20. Computing apparatus, comprising:
2 one or more volumes of data storage, arranged to
3 store data; and

4 a plurality of computing nodes, linked to access the
5 volumes of data storage using a parallel file system, and
6 arranged so as to select a first one of the nodes to
7 serve as a session manager node and to select a second
8 one of the nodes to serve as a session node for a data
9 management application, so that a session of the data
10 management application is created on the session node by
11 sending a message from the session node to the session
12 manager node, causing the session manager node to
13 distribute information regarding the session among the
14 nodes in the cluster, responsive to which the session
15 node receives events from the nodes in the cluster for
16 processing by the data management application.

1 21. Apparatus according to claim 20, wherein the session
2 manager node is arranged to store the information
3 regarding the session.

1 22. Apparatus according to claim 21, wherein the
2 information regarding the session is stored at both the
3 session node and at the session manager node, and wherein
4 following a failure at the session node, the stored

5 information is received from the session manager node in
6 order to recover the session.

1 23. Apparatus according to claim 22, wherein at least a
2 portion of the information regarding the session is
3 stored at one or more of the nodes in the cluster other
4 than at the session node and the session manager node,
5 and wherein following a failure at the first one of the
6 nodes, the nodes are arranged to select a third one of
7 the nodes to serve as the session manager node, and to
8 collect the information from at least one of the session
9 nodes and the other nodes in the cluster at which the
10 information is stored for use by the third one of the
11 nodes in serving as the session manager node.

1 24. Apparatus according to claim 20, wherein the session
2 is created in accordance with a data management
3 application programming interface (DMAPI) of the parallel
4 file system, and wherein sending the message invokes a
5 session management function of the DMAPI on the session
6 manager node.

1 25. Apparatus according to claim 20, wherein the
2 information comprises a list of the events in each file
3 system of relevance to the data management application
4 and respective dispositions of the events on the list,
5 and wherein the nodes are arranged to report to the
6 session node the events appearing on the list responsive
7 to the dispositions.

1 26. Apparatus according to claim 25, wherein the
2 dispositions can be set by any of the nodes.

1 27. Apparatus according to claim 25, wherein the session
2 is one of a plurality of sessions in the cluster, and

3 wherein the session manager node coordinates a consistent
4 partitioning of the dispositions among the sessions.

1 28. Apparatus according to claim 20, wherein the nodes
2 are arranged so that a plurality of the nodes can be
3 selected to serve as respective session nodes in a
4 plurality of data management sessions, and wherein the
5 session manager node is informed of the plurality of the
6 sessions, causing the session manager node to distribute
7 the information regarding the plurality of the sessions.

1 29. Apparatus according to claim 28, wherein the first
2 one of the nodes is arranged to serve as one of the
3 session nodes, in addition to serving as the session
4 manager node.

1 30. Apparatus according to claim 28, wherein the
2 plurality of the nodes comprises multiple session nodes
3 selected for a distributed data management application
4 running in the cluster.

1 31. Apparatus according to claim 28, wherein the
2 plurality of the nodes comprises the second one of the
3 nodes, selected to serve as the respective session node
4 for a first data management application, and a third one
5 of the nodes selected to serve as the respective session
6 node for a second data management application.

1 32. Apparatus according to claim 20, wherein the session
2 node is arranged to modify the session by sending a
3 further message to the session manager node, causing the
4 session manager node to distribute a notification
5 regarding the modified session to the nodes in the
6 cluster.

1 33. Apparatus according to claim 20, wherein the session
2 node is arranged to destroy the session by sending a
3 further message to the session manager node, causing the
4 session manager node to distribute a notification
5 regarding the destroyed session to the nodes in the
6 cluster.

1 34. Apparatus according to claim 20, wherein the data
2 management application comprises a data migration
3 application, for freeing storage space on at least one of
4 the volumes of data storage.

1 35. Apparatus according to claim 20, wherein the
2 information comprises a session identifier, generated at
3 the session manager node, which is unique within the
4 cluster.

1 36. Computing apparatus, comprising:

2 one or more volumes of data storage, arranged to
3 store data; and

4 a plurality of computing nodes, linked to access the
5 volumes of data storage using a parallel file system, and
6 arranged to initiate sessions of a parallel data
7 management application on the plurality of the nodes, so
8 that an instance of the data management application runs
9 on each of the nodes, and a data management event is
10 generated responsive to a request submitted to the
11 parallel file system on at least one of the nodes to
12 perform a file operation on a file in one of the volumes
13 of data storage, causing the event to be handled by the
14 instance of the data management application running on
15 the at least one of the nodes.

1 37. Apparatus according to claim 36, wherein the at
2 least one of the nodes is arranged to send an event
3 message to the other nodes, so as to inform the data
4 management application sessions on the other nodes of the
5 event.

1 38. Apparatus according to claim 36, wherein the data
2 management event is generated by a user application
3 running on the at least one of the nodes, which submits
4 the request.

1 39. A computer software product for use in a cluster of
2 computing nodes having shared access to one or more
3 volumes of data storage using a parallel file system, the
4 product comprising a computer-readable medium in which
5 program instructions are stored, which instructions, when
6 read by the computing nodes, cause a first one of the
7 nodes to be selected to serve as a session manager node
8 and a second one of the nodes to be selected to serve as
9 a session node for a data management application, and
10 cause a session of the data management application to be
11 created on the session node by sending a message from the
12 session node to the session manager node, causing the
13 session manager node to distribute information regarding
14 the session among the nodes in the cluster, responsive to
15 which the session node receives events from the nodes in
16 the cluster for processing by the data management
17 application.

1 40. A product according to claim 39, wherein the
2 instructions cause the session manager node to store the
3 information regarding the session.

1 41. A product according to claim 40, wherein the
2 instructions cause the information regarding the session
3 to be stored at both the session node and at the session
4 manager node, and wherein following a failure at the
5 session node, the instructions cause the stored
6 information to be received from the session manager node
7 in order to recover the session.

1 42. A product according to claim 41, wherein the
2 instructions cause at least a portion of the information
3 regarding the session to be stored at one or more of the
4 nodes in the cluster other than at the session node and
5 the session manager node, and wherein following a failure
6 at the first one of the nodes, the instructions cause the
7 nodes to select a third one of the nodes to serve as the
8 session manager node, and to collect the information from
9 at least one of the session node and the other nodes in
10 the cluster at which the information is stored for use by
11 the third one of the nodes in serving as the session
12 manager node.

1 43. A product according to claim 39, wherein the product
2 comprises a data management application programming
3 interface (DMAPI) of the parallel file system, and
4 wherein the instructions cause the session node to send
5 the message by invoking a session management function of
6 the DMAPI on the session manager node.

1 44. A product according to claim 39, wherein the
2 information comprises a list of the events in each file
3 system of relevance to the data management application
4 and respective dispositions of the events on the list,
5 and wherein the instructions cause the nodes to report to

6 the session node the events appearing on the list
7 responsive to the dispositions.

1 45. A product according to claim 44, wherein the
2 dispositions can be set by any of the nodes.

1 46. A product according to claim 44, wherein the session
2 is one of a plurality of sessions in the cluster, and
3 wherein the session manager node coordinates a consistent
4 partitioning of the dispositions among the sessions.

1 47. A product according to claim 39, wherein the
2 instructions are such as to cause a plurality of the
3 nodes to be selected to serve as respective session nodes
4 in a plurality of data management sessions, and to cause
5 the session nodes to inform the session manager node of
6 the plurality of the sessions, causing the session
7 manager node to distribute the information regarding the
8 plurality of the sessions.

1 48. A product according to claim 47, wherein the
2 instructions allow the first one of the nodes to serve as
3 one of the session nodes, in addition to serving as the
4 session manager node.

1 49. A product according to claim 47, wherein the
2 plurality of the nodes comprises multiple session nodes
3 selected for a distributed data management application
4 running in the cluster.

1 50. A product according to claim 47, wherein the
2 plurality of the nodes comprises the second one of the
3 nodes, selected to serve as the respective session node
4 for a first data management application, and a third one
5 of the nodes selected to serve as the respective session
6 node for a second data management application.

1 51. A product according to claim 39, wherein the
2 instructions cause the session node to modify the session
3 by sending a further message to the session manager node,
4 causing the session manager node to distribute a
5 notification regarding the modified session to the nodes
6 in the cluster.

1 52. A product according to claim 39, wherein the
2 instructions cause the session node to destroy the
3 session by sending a further message to the session
4 manager node, causing the session manager node to
5 distribute a notification regarding the destroyed session
6 to the nodes in the cluster.

1 53. A product according to claim 39, wherein the data
2 management application comprises a data migration
3 application, for freeing storage space on at least one of
4 the volumes of data storage.

1 54. Apparatus according to claim 39, wherein the
2 information comprises a session identifier, generated at
3 the session manager node, which is unique within the
4 cluster.

1 55. A computer software product for use in a cluster of
2 computing nodes having shared access to one or more
3 volumes of data storage using a parallel file system, the
4 product comprising a computer-readable medium in which
5 program instructions are stored, which instructions, when
6 read by the computing nodes, cause sessions of a parallel
7 data management application to be initiated on the
8 plurality of the nodes, so that an instance of the data
9 management application runs on each of the nodes, and a
10 data management event is generated responsive to a

11 request submitted to the parallel file system on at least
 12 one of the nodes to perform a file operation on a file in
 13 one of the volumes of data storage, causing the event to
 14 be handled by the instance of the data management
 15 application running on the at least one of the nodes.

1 56. A product according to claim 55, wherein the
 2 instructions cause the at least one of the nodes to send
 3 an event message to the other nodes, so as to inform the
 4 data management application sessions on the other nodes
 5 of the event.

1 57. Apparatus according to claim 55, wherein the data
 2 management event is generated by a user application
 3 running on the at least one of the nodes, which submits
 4 the request.